

Generalizing Plans to New Environments in Relational MDPs

Thomas Dickerson, Nakul Gopalan, David Meierfrankenfeld

December 12, 2013



BROWN

Generalizing plans

- learn to plan in one domain
- use this learned knowledge in the next domain using generalizations
- generalizations can be in the form of class based value functions
- this needs the classes to have the same dynamics across domains
- example: Freecraft!



BROWN

Freecraft

- state of an object is defined by its health ($0 - 4$)
- each footman can attack any one of its enemies
- an enemy can attack only its chosen footman
- attacking an object reduces its health
- with n footmen and n enemies the size of the state space is 5^{2n}
- the action space has size n^n



BROWN

Solving MDPs using linear programs

- a general MDP linear program (LP)

$$\text{Minimize: } \sum_i \alpha(s_i) V(i)$$

$$\text{Subject to: } V(i) \geq R(s_i, a) + \gamma \sum_k P(s'_k | s_i, a) V_k$$

- it has an exponential number of variables 5^{2n} and constraints $5^{2n}n^n$ for our problem
- objective function is often written as a linear combination of basis functions
- there are still exponential number of constraints



BROWN

Factored MDPs

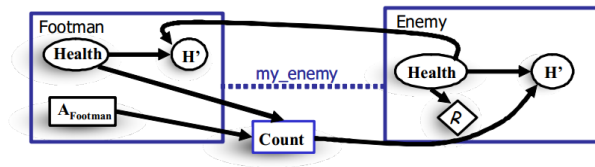
- in some problems the reward and value functions can be written as a sum of subfunctions
- these subfunctions are dependent on a small number of state variables
- in this case an LP can be constructed with the number of constraints exponential in the induced width of the dependency graph¹
- this reduces the size of the LP dramatically

¹Efficient Solution Algorithms for Factored MDPs, Guestrin et al..



Relational MDPs

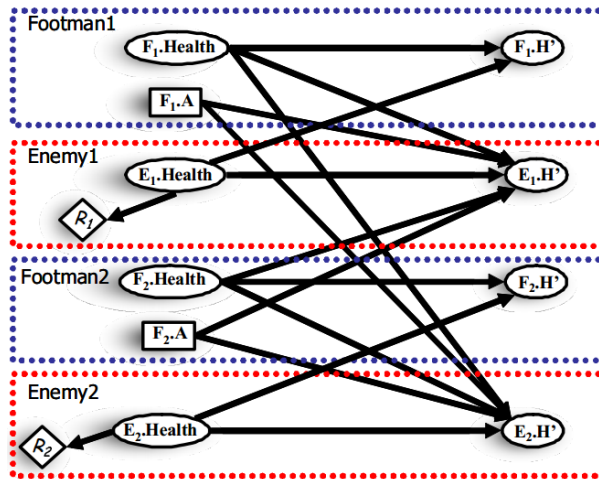
- RMDP defines how state and reward of an object depends on the states and actions of objects linked to it.
- provides class based transition dynamics
- this class structure can be used to define value functions
- class based value functions can be summed up as approximate value functions



Behaviour in tactical domain²

²Generalizing Plans to New Environments in Relational MDPs , Guestrin et al..

Corresponding factored MDPs



Factored MDPs and Value functions

- value function for tactical domain

$$\mathbf{V}_{2\text{vs}2}(F1.H, E1.H, F2.H, E2.H) = V_{F1}(F1.H, E1.H) + V_{E1}(E1.H) \\ + V_{F2}(F2.H, E2.H) + V_{E2}(E2.H)$$

- general global value function

$$\mathbf{V}(s) = \sum_{C \in \mathcal{C}} \sum_{o \in \mathcal{O}[C]} V_C(S[T_o]) ,$$

where o is an object of class C and T_o are its links



BROWN

LP to solve for Class based value functions

Variables: $V_C(t_C) : C \in \mathcal{C}, t_C \in \text{Dom}[\mathbf{T}_C]$

Minimize: $\sum_{C \in \mathcal{C}} \sum_{o \in O[C]} \sum_{t_o \in T_o} P(t_o) V_C(t_o)$

Subject to: $\sum_{C \in \mathcal{C}} \sum_{o \in O[C]} V_C(s[T_o]) \geq \sum_{o \in O[C]} R^C(s[S_o], a[o.A]) +$
 $\gamma \sum_{s'} P(s'|s, a) \sum_{C \in \mathcal{C}} \sum_{o \in O[C]} V_C(s'[T_o])$



BROWN

Implementation

- implemented the model domain on Netlogo and Burlap
- implemented joint action, class based value and reward framework in Burlap
- implemented a class based value function solver using LP in Burlap

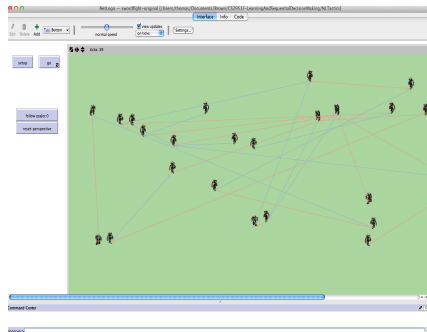


BROWN

the only plots we have



Original Plot



Our Plot



BROWN

Not really

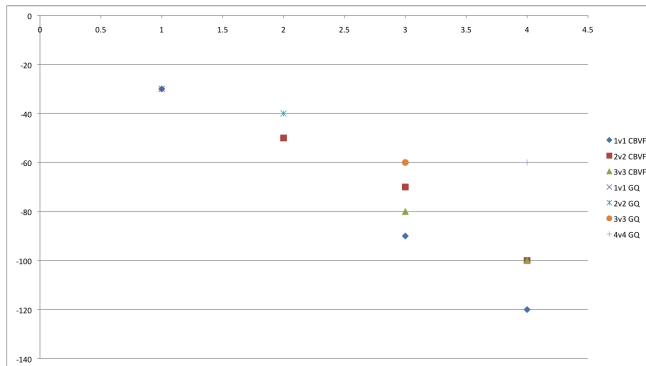
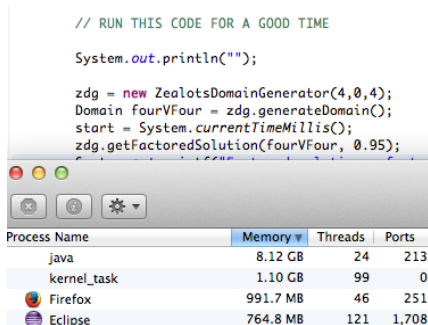
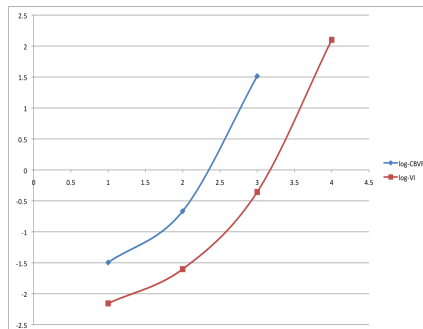


Figure: players vs total rewards

Memory



Memory usage on a fancy apple computer!



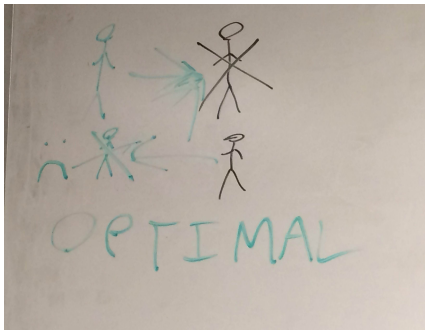
Runtime Plot log time vs # agents per side



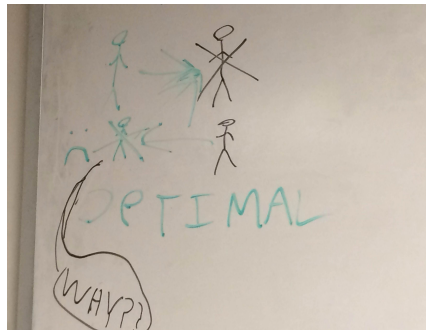
BROWN

Pain slides!

```
keys.contains(keys.iterator().next()) -> false
```



Why??



Beating the dead horse



BROWN

- The Good
 - ① class based value functions to approximate MDP as an idea
 - ② solving LPs with exponential constraints
- The Bad
 - ① tactical domains are oversimplified to favour the Footman class
 - ② the paper claims the values do not generalize well to larger set of objects
- The Ugly
 - ① lack of clarity in the paper about the LP solving process
 - ② no software released

